

日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日
Date of Application:

2003年 1月 9日

出 願 番 号
Application Number:

特願2003-002981

[ST.10/C]:

[JP 2003-002981]

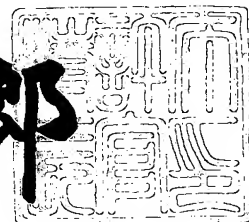
出 願 人
Applicant(s):

沖電気工業株式会社

2003年 4月 1日

特 許 庁 長 官
Commissioner,
Japan Patent Office

太田信一郎



出証番号 出証特2003-3022447

31759-190543
11/1/2003

【書類名】 特許願

【整理番号】 KN002566

【提出日】 平成15年 1月 9日

【あて先】 特許庁長官 太田 信一郎 殿

【国際特許分類】 G06F 15/40

【発明者】

【住所又は居所】 東京都港区虎ノ門1丁目7番12号 沖電気工業株式会
社内

【氏名】 池野 篤司

【特許出願人】

【識別番号】 000000295

【氏名又は名称】 沖電気工業株式会社

【代表者】 篠塚 勝正

【代理人】

【識別番号】 100090620

【弁理士】

【氏名又は名称】 工藤 宣幸

【先の出願に基づく優先権主張】

【出願番号】 特願2002-187698

【出願日】 平成14年 6月27日

【手数料の表示】

【予納台帳番号】 013664

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9006358

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 情報区分装置、方法及びプログラム、並びに、情報区分プログラムを記録した記録媒体

【特許請求の範囲】

【請求項 1】 入力された電子文書中の情報を区分する情報区分装置において、

分割行に表れ得る所定の文字列を規定する 1 又は複数の分割パターンを記憶している分割パターン記憶手段と、

入力された上記電子文書を上記分割パターン記憶手段に記憶されている上記分割パターンと照合して、上記電子文書を、複数の部分文書に分割する文書分割手段とを備える

ことを特徴とする情報区分装置。

【請求項 2】 分類を特定し得る所定の文字列を規定する、分類情報が付与されている複数のラベリングパターンを記憶しているラベリングパターン記憶手段と、

上記文書分割手段によって分割された上記各部分文書をそれぞれ、上記ラベリングパターン記憶手段に記憶されている上記ラベリングパターンと照合して、分類情報を付与するラベリング手段とをさらに備える

ことを特徴とする請求項 1 に記載の情報区分装置。

【請求項 3】 入力された上記電子文書の種類を判別する文書種類判別手段をさらに備え、

上記文書分割手段が、判別された文書種類用の上記分割パターンを用いて部分文書への分割を行うと共に、

上記ラベリング手段が、判別された文書種類用の上記ラベリングパターンを用いて分類情報の付与を行う

ことを特徴とする請求項 2 に記載の情報区分装置。

【請求項 4】 入力された上記電子文書における、同様な文字列を同様な位置に含む複数行の存在を認識して、上記分割パターンを生成し、上記分割パター

ン記憶手段に登録する分割パターン生成手段をさらに備えることを特徴とする請求項 1 ～ 3 のいずれかに記載の情報区分装置。

【請求項 5】 入力された電子文書中の情報を区分する情報区分方法において、

入力された上記電子文書を、分割行に表れ得る所定の文字列を規定する分割パターンと照合して、上記電子文書を、複数の部分文書に分割する文書分割工程を有する

ことを特徴とする情報区分方法。

【請求項 6】 上記文書分割工程によって分割された上記各部分文書をそれぞれ、分類を特定し得る所定の文字列を規定する、分類情報が付与されているラベリングパターンと照合して、分類情報を付与するラベリング工程をさらに有することを特徴とする請求項 5 に記載の情報区分方法。

【請求項 7】 入力された上記電子文書の種類を判別する文書種類判別工程をさらに有し、

上記文書分割工程が、判別された文書種類用の上記分割パターンを用いて部分文書への分割を行うと共に、

上記ラベリング工程が、判別された文書種類用の上記ラベリングパターンを用いて分類情報の付与を行う

ことを特徴とする請求項 6 に記載の情報区分方法。

【請求項 8】 入力された上記電子文書における、同様な文字列を同様な位置に含む複数行の存在を認識して、上記分割パターンを生成して登録する分割パターン生成工程をさらに有することを特徴とする請求項 5 ～ 7 のいずれかに記載の情報区分方法。

【請求項 9】 請求項 5 ～ 7 のいずれかに記載の情報区分方法の各工程をコンピュータが処理し得るコードで記述したことを特徴とする情報区分プログラム。

【請求項 10】 請求項 9 の情報区分プログラムを記録していることを特徴とする記録媒体。

【発明の詳細な説明】

【 0 0 0 1 】

【発明の属する技術分野】

本発明は、情報区分装置、方法及びプログラム、並びに、情報区分プログラムを記録した記録媒体に関し、特に、複数の情報が記載されている電子文書内の情報を分割して分類しようとするものである。

【 0 0 0 2 】

【従来の技術】

近年、インターネット等のネットワーク技術の普及により国内外の大量の電子文書へのアクセスが可能になり、大量の電子文書情報を分類する等の知的作業の自動化の必要性が高まってきている。

【 0 0 0 3 】

昨今発展を見せている電子文書の入手方法の一つに、メールマガジン（メールによる雑誌・新聞に類したもの）があげられる。これは、購読希望者に複数の情報をまとめて一つの電子メールに記載して送るというものである。

【 0 0 0 4 】

このような電子メールは、複数の情報を記載した電子文書と見なすことができ、その情報を分類するためには電子文書内の各情報を適切に分割してやる必要がある。

【 0 0 0 5 】

【特許文献 1】 特開 2 0 0 0 - 2 8 5 1 4 0 号公報

特許文献 1 には、文書データの構造情報（HTML のタグや文字のフォント情報など）を基に文書データを分割する手段や、文書要素（例えば単語）や要素付随情報（例えば品詞）を基に文書データを分割する手段を設けることにより、情報の分類の一助としている装置の例が示されている。

【 0 0 0 6 】

【発明が解決しようとする課題】

しかしながら、上記公報記載の装置では、メールマガジンのように明確な構造情報を持っていない電子文書には適用できないという問題がある。

【 0 0 0 7 】

また、仮に、あるメールマガジンを適切に分割する情報を指定したとしても、複数のメールマガジンを受け取っている場合、各々が異なる種類の分割情報（分割パターン）を必要とする可能性が高く、メールマガジンの種類によっては、適切な分割パターンを選択して分割することができないという課題がある。

【 0 0 0 8 】

さらに、受け取るメールマガジンが増加すれば、分割パターンの種類も増加するが、それらを人手で指定するのは手間がかかるという課題がある。

【 0 0 0 9 】

そのため、メールマガジン等のような明確な構造情報を持っていない電子文書の各情報を、適切に分割できる情報区分装置等が望まれている。

【 0 0 1 0 】

【課題を解決するための手段】

第 1 の本発明の情報区分装置は、入力された電子文書中の情報を区分するものであって、（１）分割行に表れ得る所定の文字列を規定する 1 又は複数の分割パターンを記憶している分割パターン記憶手段と、（２）入力された上記電子文書を上記分割パターン記憶手段に記憶されている上記分割パターンと照合して、上記電子文書を、複数の部分文書に分割する文書分割手段とを備えることを特徴とする。

【 0 0 1 1 】

第 2 の本発明の情報区分方法は、入力された電子文書中の情報を区分するものであって、入力された上記電子文書を、分割行に表れ得る所定の文字列を規定する分割パターンと照合して、上記電子文書を、複数の部分文書に分割する文書分割工程を有することを特徴とする。

【 0 0 1 2 】

第 3 の本発明の情報区分プログラムは、第 2 の本発明の情報区分方法の工程をコンピュータが処理し得るコードで記述したことを特徴とする。

【 0 0 1 3 】

第 4 の本発明の記録媒体は、第 3 の本発明の情報区分プログラムを記録していることを特徴とする。

【 0 0 1 4 】

【発明の実施の形態】

(A) 第 1 の実施形態

以下、本発明による情報区分装置、方法及びプログラム、並びに、情報区分プログラムを記録した記録媒体の第 1 の実施形態を図面を参照しながら詳述する。

【 0 0 1 5 】

(A-1) 第 1 の実施形態の構成

図 1 は、第 1 の実施形態の情報区分装置の機能的構成を示すブロック図である。例えば、第 1 の実施形態の情報区分装置は、通信機能を有するパソコン等の情報処理装置に対し、CD-ROM やフロッピー（登録商標）ディスク等の記録媒体に記録されている情報区分プログラムをインストールすることで実現されるが、機能的には、図 1 で表すことができる。

【 0 0 1 6 】

図 1 において、第 1 の実施形態の情報区分装置は、文書種類判別部 1 と、文書分割部 2 と、ラベリング部 3 と、判別パターンデータ記憶部 4 と、分割パターンデータ記憶部 5 と、ラベリングパターンデータ記憶部 6 とを有する。

【 0 0 1 7 】

文書種類判別部 1 は、判別パターンデータ記憶部 4 の判別パターンデータを参照して、適用すべき分割パターンとラベリングパターンを決定するために、入力された電子文書（適宜、文書と呼ぶ）の種類を判別するものである。

【 0 0 1 8 】

なお、この第 1 の実施形態では、複数の全く異なる情報が一つの電子文書内に含まれている電子文書（例えばニュースのメールマガジン）を入力対象としているものである。また、構造情報は持たないが、人間が簡単に認識できるように記号などの表層情報を用いて明示的に内容の区切りが記述されている電子文書を入力対象としているものである。

【 0 0 1 9 】

文書分割部 2 は、文書種類判別部 1 の判別結果（すなわち、電子文書の種類）により決定された、分割パターンデータ記憶部 5 中の分割パターンデータを適用

して、入力された電子文書を分割するものである。

【 0 0 2 0 】

ラベリング部 3 は、文書種類判別部 1 の判別結果（すなわち、電子文書の種類）結果により決定された、ラベリングパターンデータ記憶部 6 中のラベリングパターンデータを適用して、文書分割部 2 により分割された入力文書の各部分に対してラベリングを行なうものである。

【 0 0 2 1 】

判別パターンデータ記憶部 4 に記憶されている判別パターンデータは、文書種類判別部 1 が電子文書の種別を判別するためのデータの集合である。最も単純な形式の判別パターンとしては、特定の文字列（例えば、メールマガジンであれば、メールマガジンのタイトルや I D 番号）が挙げられる。

【 0 0 2 2 】

図 2 は、判別パターンデータの一例を示している。各レコードは、文書種類と、その文書種類に適用する判別パターンとを含んでいる。図 2 に示すように、ある種類の電子文書に対し、複数の判別パターンデータが存在していても構わない。

【 0 0 2 3 】

分割パターンデータ記憶部 5 に記憶されている分割パターンデータは、文書分割部 2 が電子文書を分割するためのデータであり、例えば、図 3 に示すような文書種類と分割パターンとを対応付けたデータである。図 3 の分割パターンは、正規表現で記載されているので、パターン中の記号「^」は「行頭」、「.」は「任意の一文字」、「*」は「直前の文字が 0 回以上出現する」ことを意味している。例えば、図 3 における「^====.*」は、「行頭から半角のイコール記号『=』が 4 回出現した後にある文字が 0 回以上出現する」というパターンを示していることになる。図 3 に示すように、ある種類の電子文書に対し、複数の分割パターンデータが存在していても構わない。また、電子文書の種類を問わずに適用する分割パターンデータを設けていても良い。

【 0 0 2 4 】

ラベリングパターンデータ記憶部 6 に記憶されているラベリングパターンデー

タは、文書分割部 2 が分割した電子文書の各部分（各情報）に対して、ラベリング部 3 が分類情報を付与する（ラベリングを行なう）ためのデータであり、図 4 に示すような、文書種類と、ラベリングパターンと、ラベル名とを対応付けたデータの集合である。図 4 に示すラベリングパターンも、正規表現で記載されている。図 4 に示すように、ある種類の電子文書に対し、通常、複数のラベリングパターンデータが存在する。また、電子文書の種類を問わずに適用するラベリングパターンデータを設けていても良い。

【 0 0 2 5 】

（A-2）第 1 の実施形態の動作

以下、第 1 の実施形態の情報区分装置の動作（情報区分方法）を、各構成要素 1 ～ 3 毎の動作で説明する。

【 0 0 2 6 】

まず、文書種類判別部 1 の動作を説明する。

【 0 0 2 7 】

文書種類判別部 1 は、判別パターンデータ記憶部 4 に記憶されている各パターンデータを用いて、入力された電子文書内をパターンマッチさせることにより文書種類を判別する。なお、入力文書は、ネットワークを介して取り込んでも良く、記憶媒体から取り出しても良く、その入力方法は任意である。

【 0 0 2 8 】

ここで、入力文書が図 5 に示すような電子文書であった場合には、図 2 における第 1 番目や第 2 番目のパターンデータの存在により、図 5 の電子文書は「ビジネスメールマガジン 1」という種別であると判別される。

【 0 0 2 9 】

なお、複数のパターンデータがマッチし、かつ、その判別結果が矛盾する場合には、多数決（マッチ数が多いもの）により決定したり、矛盾が生じる旨をユーザに通知するなどの機能を設けても良い。

【 0 0 3 0 】

次に、文書分割部 2 の動作を説明する。

【 0 0 3 1 】

文書分割部 2 は、上述したように、分割パターンデータ記憶部 5 に記憶されている、判別された文書種類の各分割パターンデータを用いて、入力された電子文書を複数の部分文書（情報）に分割する。

【 0 0 3 2 】

図 5 の電子文書が、文書種類判別部 1 によって「ビジネスメールマガジン 1」という種別と判別されたので、図 3 の第 1 番目及び第 2 番目の分割パターンが適用可能である。すなわち、（１）先頭から「-」（半角のハイフン）が一定数以上連続している、（２）先頭から「=」（半角の等号）が一定数以上連続している、の部分が分割パターンとなるので、その位置（行）で入力文書を部分文書（情報）に分割する。

【 0 0 3 3 】

分割後の各部分文書は、データ全般を記憶している記憶装置上に元データとは別に記憶されることになる。なお、各部分文書の記憶部は、文書分割部 2 に含まれているように、図 1 では示している。

【 0 0 3 4 】

また、分割に用いた分割パターンそのものは、（１）分割後の部分文書には含めない（分割パターンは削除される）、（２）分割位置の前後の部分文書のいずれかに含める、（３）分割位置の前後の両方の部分文書に含める（パターンは複製される）、のいずれかの方法を適用する。

【 0 0 3 5 】

分割パターンの取扱いについて（２）の方法を適用した場合には、図 5 の入力文書は、図 6 に示すような 5 個の部分文書に分割される。

【 0 0 3 6 】

次に、ラベリング部 3 の動作を説明する。

【 0 0 3 7 】

ラベリング部 3 は、上述したように、ラベリングパターンデータ記憶部 6 に記憶されている、判別された文書種類の各ラベリングパターンデータを用いて、パターンがマッチした部分文書をラベリングする。

【 0 0 3 8 】

図 5（図 6）の電子文書が文書種類判別部 1 によって「ビジネスメールマガジン 1」という種別と判別されたので、図 4 の第 1 番目～第 4 番目のラベリングパターンデータが利用され、その結果、部分文書 1 に対して「広告」、部分文書 2 に対して「タイトル」、部分文書 3 及び 4 に対して「記事本文」、部分文書 5 に対して「注釈」のようにラベリングされる。

【 0 0 3 9 】

例えば、部分文書 1 には、「――PR―」というパターンが存在するので、図 4 の第 2 番目の行が適用され、「広告」とラベリングされる。これらのラベル情報は、各部分文書と組にして保持される。

【 0 0 4 0 】

ラベル情報を有する部分文書の情報は、ユーザの操作等に応じて、表示出力されたり、印刷出力されたり、他へ送信されたりする。この際、ユーザは、例えば、記事本文だけを指定して出力させたりすることもできる。また、ラベル情報を有する部分文書の情報は、さらなる加工処理が実行されても良い。例えば、記事本文に対して要約作成処理を施すようにしても良い。

【 0 0 4 1 】

（A-3）第 1 の実施形態の効果

以上のように、第 1 の実施形態によれば、簡単なパターンによる分割パターンデータやラベリングパターンデータを用意するだけで、XML や HTML や SGML 等で記述されたような明確な構造を持つ電子文書ではなくても、電子文書を分割して分類することができる。

【 0 0 4 2 】

しかも、文書種類判別部を設けたので、複数の分割パターンを管理しておき、様々な種類の電子文書を対象に電子文書を分割して分類することができる。

【 0 0 4 3 】

（B）第 2 の実施形態

次に、本発明による情報区分装置、方法及びプログラム、並びに、情報区分プログラムを記録した記録媒体の第 2 の実施形態を図面を参照しながら詳述する。

【 0 0 4 4 】

(B-1) 第2の実施形態の構成

図7は、第2の実施形態の情報区分装置の機能的構成を示すブロック図であり、第1の実施形態に係る図1との同一、対応部分には、同一符号を付して示している。

【0045】

第2の実施形態の情報区分装置は、第1の実施形態の構成に、分割パターン生成部7を付加した構成となっている。

【0046】

分割パターン生成部7は、入力された電子文書を基に分割パターンを生成するものである。分割パターン生成部7によって生成された分割パターンは、文書分類判別部1によって判別された文書種類に対応付けられ、分割パターンデータとして分割パターンデータ記憶部5に記憶される。

【0047】

分割パターン生成部7以外の部分は、第1の実施形態と同様の機能を担っているので、その説明は省略する。

【0048】

(B-2) 第2の実施形態の動作

第1の実施形態と動作が異なるのは分割パターン生成部7の動作だけなので、以下では、その動作のみを、図8のフローチャートを参照しながら説明する。

【0049】

分割パターン生成部7は、入力文書が与えられると、入力文書を行ごとに分割する(ステップ801)。次に、先頭から所定番目(例えば30番目)の文字の全てが一致する行のグループを作ると共に、その行グループに属する行数も計数しておく(ステップ802)。

【0050】

例えば、上述した図5の電子文書が入力文書である場合、ステップ802の処理を終えた段階では、図9に示すような行グループが作成される。

【0051】

その後、分割パターン生成部7は、複数のメンバ(行)(ここでは2以上とす

る)を持つ行グループのみを選択してパターン記述を行う(ステップ803)。最も簡単なパターン記述法は文字列そのものであるが、必要に応じて正規表現などに書き改めるなどの手法を用いても良く、文書分割部2が理解できる形式を出力するものであれば特に手法は問わない。

【0052】

その後、分割パターン生成部7は、文書種類判別部1から、文書種類のデータを取り込んで分割パターンデータを完成させて分割パターンデータ記憶部5に登録する(ステップ804)。なお、文書種類のデータを含まない分割パターンデータを登録するようにしても良い。

【0053】

上述したステップ802の行一致を判別するための文字数や、ステップ803の登録に値するかを判別するためのメンバ(行)数は自由に設定しても良い。また、ステップ802において「先頭から複数文字」としているが、「末尾から」であっても良く、「先頭および末尾から」であっても良く、「先頭や末尾に関係なく」であっても良い。また、それらを自由に設定できる形式であっても良い。

【0054】

(B-3) 第2の実施形態の効果

第2の実施形態によれば、第1の実施形態と同様な効果を奏すると共に、さらに、自動的に分割パターンデータを生成して登録することができるという効果をも奏する。

【0055】

(C) 他の実施形態

上記各実施形態においては、入力文書の分割を行った後に、各部分文書に対するラベリングを行うものを示したが、入力文書の分割及び分割された各部分文書に対するラベリングを並行して同時に行なっても良い。

【0056】

また、分割パターンデータをラベリングパターンデータの一部として用いるようにしても良い。

【0057】

上記各実施形態は、入力文書が横書き文書であるものを示したが、縦書き文書に対応できるようにしても良い。この場合、縦方向の行パターンを利用して、上記各実施形態と同様な処理を行うようにすれば良い。

【 0 0 5 8 】

また、上記各実施形態では、文書種類判別部が入力文書の種類を自動判別するものを示したが、ユーザ等が入力文書の種類を入力するものであっても良い。また、全ての分割パターンやラベリングパターンを、文書種類に関係なく、登録しておき、入力文書の種類を特定することなく、部分文書への分割、及び、分割された部分文書へのラベリングを行うようにしても良い。さらに、ある種類の入力文書専用の情報区分装置として装置を構成しても良い。

【 0 0 5 9 】

さらに、上記各実施形態の分割パターンは、その行が分割行であることを確定するものであったが、ある分割パターン（様子見分割パターン）に一致する行より所定行以内に、他の分割パターンに一致する行がないことを判明した場合に、分割行と確定するような分割パターン（様子見分割パターン）を設けるようにしても良い。

【 0 0 6 0 】

【発明の効果】

以上のように、本発明によれば、メールマガジン等のような明確な構造情報を持っていない電子文書の各情報を、適切に分割することができる。

【図面の簡単な説明】

【図 1】

第 1 の実施形態の情報区分装置の機能的構成を示すブロック図である。

【図 2】

第 1 の実施形態の判別パターンデータ例を示す説明図である。

【図 3】

第 1 の実施形態の分割パターンデータ例を示す説明図である。

【図 4】

第 1 の実施形態のラベリングパターンデータ例を示す説明図である。

【図 5】

第 1 の実施形態の動作説明に適用する入力文書例を示す説明図である。

【図 6】

図 5 の入力文書に対する文書分割処理後のデータを示す説明図である。

【図 7】

第 2 の実施形態の情報区分装置の機能的構成を示すブロック図である。

【図 8】

第 2 の実施形態の分割パターン生成部の動作を示すフローチャートである。

【図 9】

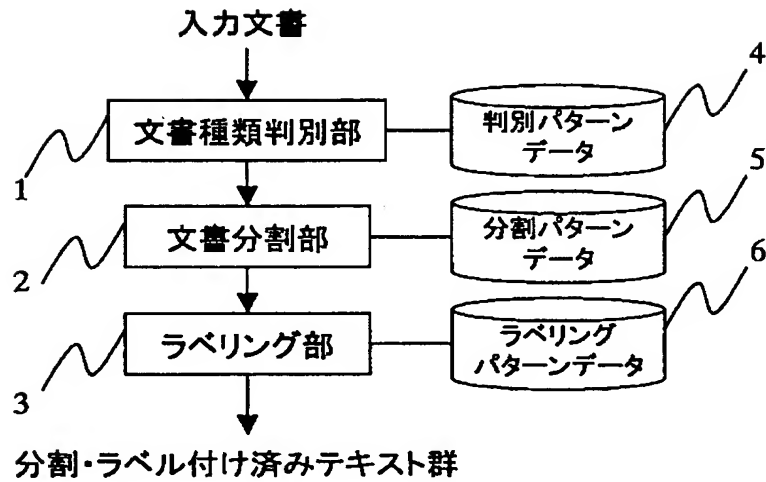
第 2 の実施形態の分割パターン生成時における入力文字のグループ化の説明図である。

【符号の説明】

1 …文書種類判別部、 2 …文書分割部、 3 …ラベリング部、 4 …判別パターンデータ記憶部、 5 …分割パターンデータ記憶部、 6 …ラベリングパターンデータ記憶部、 7 …分割パターン生成部。

【書類名】 図面

【図 1】



【図 2】

文書種類	判別パターン
ビジネスメールマガジン1	“DEFビジネスメールマガジン”
ビジネスメールマガジン1	ID=0000111
ビジネスメールマガジン2	“XYZニュース”
ビジネスメールマガジン2	MailMagazineID=999

【図 3】

文書種類	分割パターン
ビジネスメールマガジン1	^====.*
ビジネスメールマガジン1	^----.*
ビジネスメールマガジン2*
ビジネスメールマガジン2	^=-=.*

【図 4】

文書種類	ラベリングパターン	ラベル
ビジネスメールマガジン1	^●.*	記事本文
ビジネスメールマガジン1	---PR-	広告
ビジネスメールマガジン1	ID=000.*	タイトル
ビジネスメールマガジン1	^Copyright.*	注釈
ビジネスメールマガジン2	^◆.*	記事本文
ビジネスメールマガジン2	^©.*	注釈

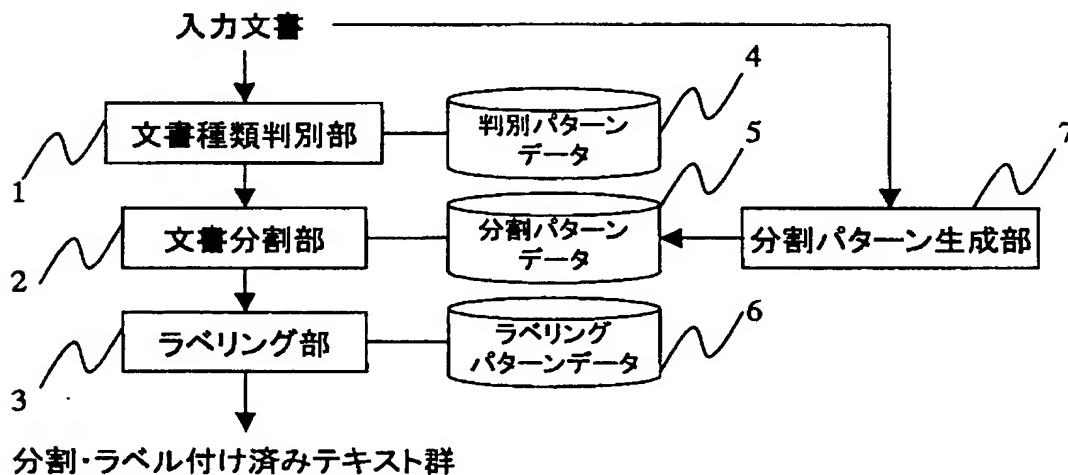
【図 5】

<p>---PR-----</p> <p>ABC Computer 新製品発売中！！</p> <p>-----</p> <p>DEF ビジネスメールマガジン 朝刊</p> <p style="text-align: right;">ID-0000111</p> <p>-----</p> <p>●GHI株式会社、JKFコーポレーションとの提携を発表</p> <p>http://www.aaa.bbb/ccc/ddd1234.html</p> <p>GHI株式会社は××日、JKFコーポレーションと**分野での提携に 合意したと発表した。...</p> <p>=====</p> <p>●株式会社LMN、第1四半期の売上高は〇〇億円</p> <p>http://www.aaa.bbb/ccc/ddd5678.html</p> <p>LMNは△△の好調な売れ行きを反映して、第1四半期の売上高は 〇〇億円と前期比**%増を...</p> <p>=====</p> <p>Copyright 2002 OPQ社 無断転載を禁ず</p>
--

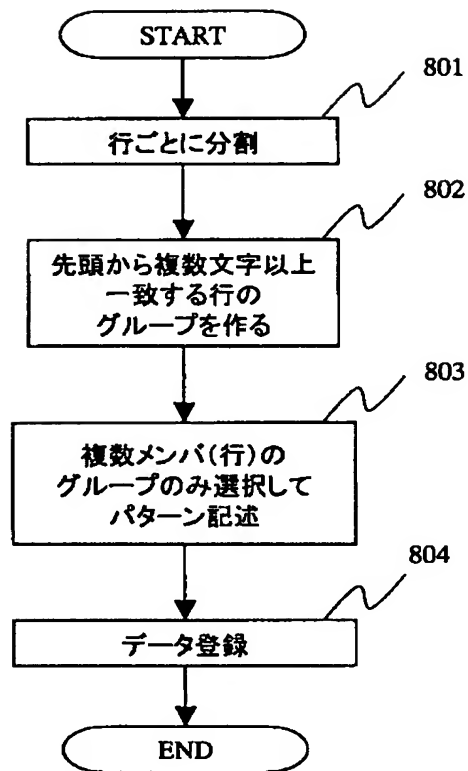
【図 6】

部分文書1	<div> <div>---PR---</div> <div>ABC Computer 新製品発売中！！</div> </div>
部分文書2	<div> <div>DEF ビジネスメールマガジン 朝刊</div> <div>ID-0000111</div> </div>
部分文書3	<div> <div>●GHI株式会社、JKFコーポレーションとの提携を発表</div> <div>http://www.aaa.bbb/ccd/ddd1234.html</div> <div>GHI株式会社は××日、JKFコーポレーションと**分野での提携に合意したと発表した。...</div> </div>
部分文書4	<div> <div>●株式会社LMN、第1四半期の売上高は〇〇億円</div> <div>http://www.aaa.bbb/ccd/ddd5678.html</div> <div>LMNは△△の好調な売れ行きを反映して、第1四半期の売上高は〇〇億円と前期比**%増を...</div> </div>
部分文書5	<div> <div>Copyright 2002 OPQ社 無断転載を禁ず</div> </div>

【図 7】



【図 8】



【図 9】

先頭からの文字	行数
----(以下30個目まで)	2
==== (以下30個目まで)	2
---PR--- (以下30個目まで)	1
...	...
http://www.aaa.bbb/ccc/ddd1234	1
...	...

【書類名】 要約書

【要約】

【課題】 メールマガジン等のような明確な構造情報を持っていない電子文書の各情報を、適切に分割する。

【解決手段】 本発明では、入力された電子文書を、分割行に表れ得る所定の文字列を規定する分割パターンと照合して、複数の部分文書に分割する。その後、分割された各部分文書をそれぞれ、分類を特定し得る所定の文字列を規定する、分類情報が付与されているラベリングパターンと照合して、分類情報を付与することが好ましい。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号 [0 0 0 0 0 0 2 9 5]

1. 変更年月日	1 9 9 0 年 8 月 2 2 日
[変更理由]	新規登録
住 所	東京都港区虎ノ門1丁目7番12号
氏 名	沖電気工業株式会社